

1 问题描述

关联规则分析是数据挖掘中活跃的研究方法之一，其目的是在一个数据集中找出各项之间的关联关系（这种关系一般没有直接在数据中表示出来）。Apriori 是关联规则分析中最常用也是最经典的挖掘频繁项集的算法，在本次作业中，我将实现 Apriori 算法，从交易数据集中发现频繁项集，并生成相应的关联规则。

2 解决方案¹

2.1 数据集的准备

为验证 Apriori 算法实现的正确性，我根据一定的规则生成了一个交易数据集，数据集的每个实例代表一条交易记录，共 100 个实例。我将商品分为食品（bread、milk、apple、orange、beer）电器（TV、PC、phone、fridge、ele_oven）和工具（scissors、stapler、plate、knife、glue）三类，各按照一定的出现及组合概率生成相应的商品交易记录，具体的参数设定可参见附录 A.1。

我认为对于交易情况的分析，仅使用模拟生成的数据集难以取得真实的结果，因此我额外在 Groceries 数据集上执行了 Apriori 算法。Groceries 数据集是内置于 R 语言的关联分析数据集，来源于某杂货店一个月的真实交易记录，包含 9835 条交易记录及 169 种商品。我将其从 R 语言包中提取出并重组为.csv 格式，再使用 Apriori 算法进行分析。

2.2 Apriori 算法

Apriori 算法是最经典的挖掘频繁项集的算法，实现了在大数据集上可行的关联规则提取，其核心思想是通过连接产生候选项与其支持度，然后通过剪枝生成频繁项集，步骤主要为：

1. 找出所有的频繁项集（支持度大于等于给定的阈值）；
2. 由频繁项集产生强关联规则（经过上个步骤后满足给定的置信度阈值的规则）。

为验证 Apriori 算法实现的正确性，我先使用蛮力算法在模拟交易数据集上运行一次，将结果与 Apriori 算法的结果进行比较。验证算法的正确性后，在 Groceries 数据集上我则直接使用 Apriori 算法进行分析。

3 实验及结果

3.1 模拟数据集

我将支持度和置信度的阈值分别设置为 0.1 及 0.6，蛮力算法的运行结果根据支持度和置信度排序后为：

¹本次作业的主要代码实现可参见附录 A.2

——频繁项集——

```
('scissors', 'stapler') , 0.15
('plate', 'stapler') , 0.15
('knife', 'stapler') , 0.15
('scissors', 'knife') , 0.15
('stapler', 'glue') , 0.16
('ele_oven', 'TV') , 0.17
('milk',) , 0.18
('apple',) , 0.18
('bread',) , 0.19
('orange',) , 0.21
('beer',) , 0.22
('PC',) , 0.22
('phone',) , 0.22
('glue',) , 0.23
('plate',) , 0.23
('scissors',) , 0.25
('fridge',) , 0.25
('ele_oven',) , 0.26
('knife',) , 0.27
('TV',) , 0.29
('stapler',) , 0.30
```

——关联规则——

```
('scissors',) --> ('stapler',) , 0.60
('scissors',) --> ('knife',) , 0.60
('plate',) --> ('stapler',) , 0.65
('ele_oven',) --> ('TV',) , 0.65
('glue',) --> ('stapler',) , 0.70
```

得到以上的参考结果后，我使用 Apriori 算法和同样的参数对模拟交易数据集进行分析，所得结果见表 1 和表 2。

表 1: Apriori 算法在模拟交易数据集下发现的频繁项集（按支持度排序）

频繁项集	支持度	频繁项集	支持度
stapler, scissors	0.15	PC	0.22
knife, scissors	0.15	phone	0.22
stapler, knife	0.15	glue	0.23
plate, stapler	0.15	plate	0.23
stapler, glue	0.16	scissors	0.25
ele_oven, TV	0.17	fridge	0.25
milk	0.18	ele_oven	0.26
apple	0.18	knife	0.27
bread	0.19	TV	0.29
orange	0.21	stapler	0.30
beer	0.22		

表 2: Apriori 算法在模拟交易数据集下发现的关联规则（按置信度排序）

关联规则	置信度
scissors \rightarrow stapler	0.60
scissors \rightarrow knife	0.60
plate \rightarrow stapler	0.65
ele_oven \rightarrow TV	0.65
glue \rightarrow stapler	0.70

对比蛮力算法和 Apriori 算法排序后的结果，容易发现二者一致，可以证明我实现的 Apriori 算法的正确性。

3.2 Groceries 数据集

确认 Apriori 算法实现的正确性后，我将其应用到真实数据上。考虑到 Groceries 数据集商品种类相对于实例较少的特点，我选择了较小的支持度（0.01），置信度选择为 0.2。运行结果参见表 3 和表 4。

表 3: Apriori 算法在 Groceries 数据集下发现的频繁项集（按支持度排序）

频繁项集	支持度	频繁项集	支持度
stapler, scissors	0.15	PC	0.22
knife, scissors	0.15	phone	0.22
stapler, knife	0.15	glue	0.23
plate, stapler	0.15	plate	0.23
stapler, glue	0.16	scissors	0.25
ele_oven, TV	0.17	fridge	0.25
milk	0.18	ele_oven	0.26
apple	0.18	knife	0.27
bread	0.19	TV	0.29
orange	0.21	stapler	0.30
beer	0.22		

表 4: Apriori 算法在 Groceries 数据集下发现的关联规则（按置信度排序）

关联规则	置信度
scissors \rightarrow stapler	0.60
scissors \rightarrow knife	0.60
plate \rightarrow stapler	0.65
ele_oven \rightarrow TV	0.65
glue \rightarrow stapler	0.70

4 结论

A 附录

A.1 模拟交易数据集的详细信息

我将模拟交易数据集的交易记录按照交易商品数分为4类，即2、3、4、5件。不同的商品件数按照不同的比例随机混合三类商品，具体混合规则可参见表5。

表 5: 模拟交易数据集生成交易记录的混合规则

商品数	混合规则（括号中数字代表各类商品在记录中所占数量）
2	(2); (1, 1)
3	(3); (2, 1)
4	(4); (3, 1); (2, 1, 1)
5	(5); (4, 1); (3, 2); (2, 2, 1)

A.2 Apriori 算法实现代码

```
1 from sklearn.cross_validation import StratifiedKFold
2 from sklearn.datasets import make_classification
3 from sklearn.cross_validation import train_test_split
4 from sklearn.metrics import classification_report
5
6 from sklearn.tree import export_graphviz
7 from sklearn.tree import DecisionTreeClassifier
8 from sklearn.ensemble import BaggingClassifier
9 from sklearn.ensemble import AdaBoostClassifier
10 from sklearn.grid_search import GridSearchCV
11
12 import numpy as np
13 import matplotlib.pyplot as plt
14 plt.rcParams['font.sans-serif']=['SimHei']
15 plt.rcParams['axes.unicode_minus']=False
16
17 ##### making dataset #####
18
19 n = 500          # number of instances
20 n_f = 30         # number of features
21 n_c = 3          # number of classes
22 inf_f = int(0.6 * n_f) # 60% real features
23 red_f = int(0.1 * n_f) # 10% redundant features
24 rep_f = int(0.1 * n_f) # 10% repeated features
25 random_seed = 1   # random seed for the experiments
26
27 X, Y = make_classification(n_samples=n, n_classes=n_c, flip_y=0.03,
28                           n_features=n_f, n_informative=inf_f, n_redundant=
29                           red_f,
30                           n_repeated=rep_f, random_state=random_seed)
```

```

30 X_train, X_test, Y_train, Y_test = \
31     train_test_split(X, Y, test_size=0.8, random_state=random_seed)
32
33 ##### making dataset #####
34
35 # directly train the model with training set & testing set
36 def exp_plain_train():
37     model = DecisionTreeClassifier(random_state=random_seed)
38     model.fit(X_train, Y_train)
39
40     pred = model.predict(X_test)
41     print(classification_report(Y_test, pred))
42
43     score = model.score(X_test, Y_test)
44     print('plain_train_score_in_testing_set:', score)
45     score = model.score(X_train, Y_train)
46     print('plain_train_score_in_training_set:', score)
47
48 # test model with cross-validation
49 def exp_cv(folds=10):
50     kfolds = StratifiedKFold(Y, n_folds=folds, random_state=
51         random_seed)
52     model = DecisionTreeClassifier(random_state=random_seed)
53     scores_train = []
54     scores_test = []
55     for train, test in kfolds:
56         model.fit(X[train], Y[train])
57         # pred = model.predict(X[test])
58
59         score = model.score(X[test], Y[test])
60         scores_test.append(score)
61         score = model.score(X[train], Y[train])
62         scores_train.append(score)
63
64     mean_test = np.array(scores_test).mean()
65     mean_train = np.array(scores_train).mean()
66     print('avg_score_with_cv_folds_%d_in_testing_set:'
67         % folds, mean_test)
68     print('avg_score_with_cv_folds_%d_in_training_set:'
69         % folds, mean_train)
70
71     return mean_test, mean_train
72
73 def plot_cv():
74     mean_tests = []
75     mean_trains = []
76     folds_sum = 11
77     for i in range(2, folds_sum):

```

```

77     m1, m2 = exp_cv(i)
78     mean_tests.append(m1)
79     mean_trains.append(m2)
80
81     # start to plot
82     x = np.arange(2, folds_sum)
83     total_width, n = 0.8, 2
84     width = total_width / n
85     x = x - (total_width - width) / 2
86
87     plt.bar(x, mean_tests, width=width,
88            facecolor='#9999ff', edgecolor='white', label=u'测试集')
89     plt.bar(x + width, mean_trains, width=width,
90            facecolor='#ffa07a', edgecolor='white', label=u'训练集')
91     for x, y1, y2 in zip(x, mean_tests, mean_trains):
92         plt.text(x - 0.05, y1 + 0.01, '%.2f' % y1, ha='center', va='
            bottom')
93         plt.text(x+width - 0.05, y2 + 0.01, '%.2f' % y2, ha='center',
            va='bottom')
94
95     plt.xlabel(u'折数')
96     plt.ylabel(u'平均_f1-score')
97     plt.ylim((0, 1.3))
98     plt.legend()
99     plt.savefig('report/img/cv_bar')
100    plt.show()
101
102    # grid search
103    def exp_grid_search(folds=10):
104        model = DecisionTreeClassifier(random_state=random_seed)
105        param_grid = {'criterion': ['gini', 'entropy'],
106                     'max_features': ['sqrt', 'log2', None],
107                     'max_depth': list(range(3, 15)),
108                     'presort': [True, False],
109                     'splitter': ['best', 'random']
110                     }
111        grid = GridSearchCV(model, param_grid, cv=folds, scoring='
            f1_weighted')
112        grid.fit(X, Y)
113
114        print(grid.best_params_)
115        print(grid.best_score_)
116
117        export_graphviz(grid.best_estimator_, filled=True, out_file='
            report/img/gs.dot')
118
119    # bagging alg
120    def bagging(cv=True):

```

```

121 bagging = BaggingClassifier(
122     DecisionTreeClassifier(random_state=random_seed),
123     n_estimators=5,                # number of models
124     random_state=random_seed,
125     bootstrap=True,
126     max_samples=1.0,              # Bootstrap sample size radio
127     bootstrap_features=True,
128     max_features=1.0,            # Bootstrap feature usage radio
129 )
130 if cv:                            # using cross-validation
131     scores_train = []
132     scores_test = []
133     kfolds = StratifiedKFold(Y, n_folds=10, random_state=
134                             random_seed)
135     for train, test in kfolds:
136         bagging.fit(X[train], Y[train])
137
138         score = bagging.score(X[test], Y[test])
139         scores_test.append(score)
140         score = bagging.score(X[train], Y[train])
141         scores_train.append(score)
142
143     mean_test = np.array(scores_test).mean()
144     mean_train = np.array(scores_train).mean()
145     print('avg_score_with_cv_folds_10_in_testing_set:', mean_test)
146     print('avg_score_with_cv_folds_10_in_training_set:',
147           mean_train)
148 else:                            # without cross-validation
149     bagging.fit(X_train, Y_train)
150     pred = bagging.predict(X_test)
151     print(classification_report(Y_test, pred))
152
153 # check the features extracted by each model
154 plt.figure(figsize=(7, 5))
155 f_n = 30
156 x = list(range(1, f_n + 1))
157 for i, f in enumerate(bagging.estimators_features_):
158     print('model%d' % (i + 1), f)
159     plt.scatter(x, f, label=u'子模型%d' % (i + 1))
160 plt.xlabel(u'特征编号')
161 plt.xticks(list(range(0, 41, 5)))
162 plt.ylabel(u'特征数值')
163 plt.legend(loc=1)
164 plt.savefig('report/img/bagging_feature%d' % len(bagging.
165           estimators_features_))
166 plt.show()

```



```

165 def plot_bagging():
166     # results
167     x = list(range(10, 101, 10))
168     y = [0.694, 0.720, 0.726, 0.730, 0.730, 0.738, 0.738, 0.738,
169         0.738, 0.744]
170     y_b1 = [0.626, 0.672, 0.702, 0.716, 0.704, 0.727, 0.735, 0.732,
171         0.736, 0.738]
172     y_b2 = [0.668, 0.711, 0.738, 0.734, 0.744, 0.746, 0.738, 0.730,
173         0.740, 0.738]
174
175     # plotting code
176     plt.figure(figsize=(6, 4))
177     ax = plt.gca()
178     ax.plot(x, y_b1, color='#90EE90', linewidth=1.7, label=u'70%特征')
179     ax.plot(x, y_b2, color='#ffa07a', linewidth=1.7, label=u'90%特征')
180     ax.plot(x, y, color='#9999ff', linewidth=1.7, label=u'100%特征')
181     ax.scatter(x, y, s=13, c='#9999ff')
182     ax.scatter(x, y_b1, s=13, c='#90EE90')
183     ax.scatter(x, y_b2, s=13, c='#ffa07a')
184     ax.grid(color='b', alpha=0.5, linestyle='dashed', linewidth=0.5)
185     plt.xlim((5, 105))
186     plt.xticks(x)
187     plt.xlabel(u'子模型数量')
188     plt.ylabel(u'平均_f1-score')
189     plt.legend()
190     plt.savefig('report/img/bagging_kline')
191     plt.show()
192
193 # boosting alg
194 def boosting(cv=True):
195     boosting = AdaBoostClassifier(
196         DecisionTreeClassifier(max_depth=3, min_samples_leaf=2,
197                                random_state=random_seed),
198         n_estimators=15, # number of models
199         algorithm='SAMME', # Advanced-Boosting
200         random_state=random_seed
201     )
202     if cv: # using cross-validation
203         scores_train = []
204         scores_test = []
205         kfolds = StratifiedKFold(Y, n_folds=10, random_state=
206             random_seed)
207         for train, test in kfolds:
208             boosting.fit(X[train], Y[train])
209
210             score = boosting.score(X[test], Y[test])
211             scores_test.append(score)
212             score = boosting.score(X[train], Y[train])

```

```

208         scores_train.append(score)
209
210     mean_test = np.array(scores_test).mean()
211     mean_train = np.array(scores_train).mean()
212     print('avg_score_with_cv_folds_10_in_testing_set:', mean_test
213           )
214     print('avg_score_with_cv_folds_10_in_training_set:',
215           mean_train)
216 else:      # without cross-validation
217     boosting.fit(X_train, Y_train)
218
219     pred = boosting.predict(X_train)
220     print(classification_report(Y_train, pred))
221
222     pred = boosting.predict(X_test)
223     print(classification_report(Y_test, pred))
224
225 # plot the relation between weights and error
226 plt.figure()
227 plt.xlabel(u'子模型权重')
228 plt.ylabel(u'错误率')
229 plt.plot(boosting.estimator_weights_, boosting.estimator_errors_)
230 plt.savefig('report/img/boosting-weight-error-%d' % len(boosting.
231               estimator_weights_))
232 plt.show()
233
234 def plot_boosting():
235     # results
236     x = list(range(10, 101, 10))
237     y1 = [0.621, 0.649, 0.649, 0.666, 0.666, 0.672, 0.672, 0.688,
238           0.688, 0.710]
239     y2 = [0.608, 0.570, 0.596, 0.616, 0.635, 0.633, 0.629, 0.650,
240           0.664, 0.668]
241     y1_tr = [0.860, 0.958, 0.984, 0.995, 0.996, 0.999, 1.000, 1.000,
242              1.000, 1.000]
243     y2_tr = [0.904, 0.936, 0.960, 0.975, 0.988, 0.993, 0.993, 0.995,
244              0.997, 0.998]
245
246     # plot testing results
247     plt.figure(figsize=(6, 4))
248     ax = plt.gca()
249     ax.plot(x, y1, color='#9999ff', linewidth=1.7, label='SAMME.R')
250     ax.plot(x, y2, color='#90EE90', linewidth=1.7, label='SAMME')
251     ax.scatter(x, y1, s=13, c='#9999ff')
252     ax.scatter(x, y2, s=13, c='#90EE90')
253     ax.grid(color='b', alpha=0.5, linestyle='dashed', linewidth=0.5)
254     plt.xlim((5, 105))
255     plt.xticks(x)

```

```

249 plt.xlabel(u'子模型数量')
250 plt.ylabel(u'平均_f1-score')
251 plt.legend()
252 plt.savefig('report/img/boosting_kline_test')
253 plt.show()
254
255 # plot training results
256 plt.figure(figsize=(6, 4))
257 ax = plt.gca()
258 ax.plot(x, y1_tr, color='#9999ff', linewidth=1.7, label='SAMME.R')
259 ax.plot(x, y2_tr, color='#90EE90', linewidth=1.7, label='SAMME')
260 ax.scatter(x, y1_tr, s=13, c='#9999ff')
261 ax.scatter(x, y2_tr, s=13, c='#90EE90')
262 ax.grid(color='b', alpha=0.5, linestyle='dashed', linewidth=0.5)
263 plt.xlim((5, 105))
264 plt.xticks(x)
265 plt.xlabel(u'子模型数量')
266 plt.ylabel(u'平均_f1-score')
267 plt.legend()
268 plt.savefig('report/img/boosting_kline_train')
269 plt.show()
270
271 if __name__ == '__main__':
272     # exp_plain_train()
273     # exp_cv()
274     # plot_cv()
275     # exp_grid_search()
276     boosting()
277     # plot_boosting()
278     # bagging()
279     # plot_bagging()

```